

Jacobian 神经网络算法

甄景贤 (King-Yin Yan)

General.Intelligence@Gmail.com

经典的神经网络 Back Prop 学习算法，它是一个 error-driven 算法，但在很多人工智能的实际应用中，不存在唯一的「理想答案」，而是根据正或负的奖励 (reward) 学习。当答案正确时，奖励 > 0 , error = 0；当答案不正确时，奖励 < 0 ，但 error 仍是不知道的（因为不知道理想答案）。简言之，就是不能用 error-driven 学习。

所以我想出了一个 reward-driven 的学习法：假设神经网络将 $\mathbf{x}_0 \mapsto \mathbf{y}_0$ ，它通常也会将 \mathbf{x}_0 的邻域 map 到 \mathbf{y}_0 的邻域。如果我们想「加强」这个映射，可以将「更大的 \mathbf{x}_0 的邻域」映射到「接近 \mathbf{y}_0 的邻域」。

这种算法对人工智能应该很重要，暂时我还想不出有什么其他办法，可以做到 [深度] 神经网络的 reward-driven 学习。

将这思想更准确化，可以将 feed-forward 神经网络的构造看成是这样的：

$$\mathbf{y} = \mathbf{F}(\mathbf{x}) \quad (1)$$

$$\mathbf{y} = \bigcirc^L \tilde{W} \dots \bigcirc^\ell \tilde{W} \dots \bigcirc^1 \tilde{W} \mathbf{x} \quad (2)$$

其中 W 代表每一层 (layer) ℓ 的矩阵。

（下面会用到） F 的反方向是：

$$\mathbf{x} = \mathbf{F}^{-1}(\mathbf{y}) \quad (3)$$

$$\mathbf{x} = \overset{1}{\rhd} \tilde{W} \dots \overset{\ell}{\rhd} \tilde{W} \dots \overset{L}{\rhd} \tilde{W} \mathbf{y} \quad (4)$$

注意： $\rhd = W^{-1}$ ， $\bigcirc = \bigcirc^{-1}$ ，形状不同。

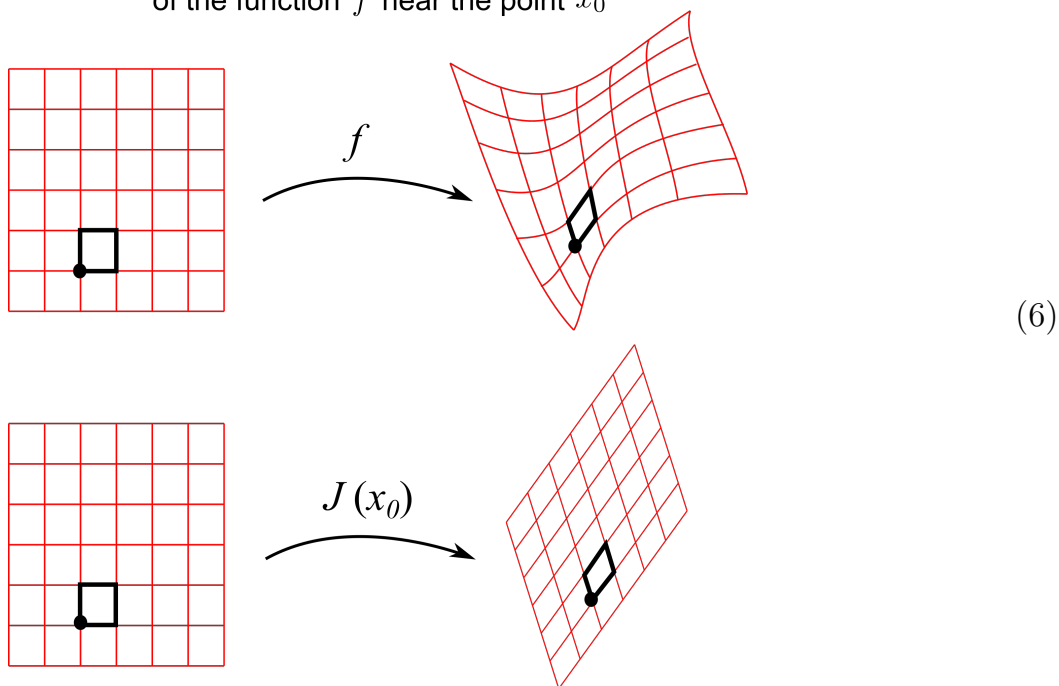
假设在 \mathbf{x} 空间有体积元 U ，经过 \mathbf{F} 变换成 \mathbf{y} 空间的体积元 V ，那么：

$$U = |J| \cdot V \quad (5)$$

$J = \left[\frac{\partial \mathbf{F}(x)}{\partial x} \right]$ 叫 Jacobian 矩阵。

在我们的情况下， $|J| = \left| \frac{\partial \mathbf{F}^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right|$ 在 \mathbf{y}_0 的值，代表「单位体积元由 $\mathbf{y}_0 \mapsto \mathbf{x}_0$ 的变化率」。
 (下面会看到， \mathbf{F} 和 \mathbf{F}^{-1} 的正/反方向不太重要，因为基本上不影响计算复杂度。)

the Jacobian $J(x_0)$ is a local linearization
 of the function f near the point x_0



每次得到**正奖励**，我们会令 Jacobian $|J|$ 增加一点：

$$|J| := \det \left[\frac{\partial \mathbf{F}^{-1}(\mathbf{y})}{\partial \mathbf{y}} \right]_{n \times n} \quad (7)$$

下标表示那是一个 $n \times n$ 矩阵。

其实 Jacobian 矩阵的意义就是：

$$J = \left[\frac{\partial \text{输出}}{\partial \text{输入}} \right] \quad (8)$$

神经网络的输入和输出都是 $\dim n$ ，所以 Jacobian 很自然是 $n \times n$ 矩阵。

用**梯度下降法**，我们需要计算这些梯度： $\left[\frac{\partial |J|}{\partial \mathbf{W}} \right]$ ，总数是网络中的 weights 的个数 $= \sum m_\ell$ 。

要用到 determinant 的微分公式：

$$\frac{d}{dt} |A(t)| = \text{tr}(\text{adj}(A) \cdot \frac{dA(t)}{dt}) \quad (9)$$

$$\text{adj}(A) := |A| \cdot A^{-1} \quad (10)$$

换句话说，对於每个权重 $w := W_{ij}^\ell$ ，我们要计算：

$$\frac{\partial}{\partial w} |J| = \text{tr}(|J| \cdot J^{-1} \cdot \left[\frac{\partial J}{\partial w} \right]) \quad (11)$$

注意： $|J|$ 和 J^{-1} 是 \mathbf{y}_0 的函数，只需在大 loop 外一次过计算。

问题是，计算 $\left[\frac{\partial J}{\partial w} \right]_{n \times n}$ 的时候：

$$\frac{\partial J}{\partial w} = \frac{\partial}{\partial w} \frac{\partial \mathbf{F}^{-1}}{\partial y} = \frac{\partial}{\partial w} \frac{\partial}{\partial \mathbf{y}} \stackrel{1}{\geq} \mathcal{O} \dots \stackrel{\ell}{\geq} \mathcal{O} \dots \stackrel{L}{\geq} \mathcal{O} \mathbf{y} \quad (12)$$

这牵涉到用 w 对 W^{-1} 的分量微分，可以想像就算计了出来也会是极复杂的。解决办法是，索性「本末倒置」，用 \geq 来定义神经网络，然后在 forward propagation 时才用 $W = \geq^{-1}$ 计算。

J 的分量写出来是：

$$J_{ij} = \frac{\partial \mathbf{F}_i^{-1}}{\partial y_j} = \frac{\partial}{\partial y_j} \left[\stackrel{1}{\geq} \mathcal{O} \dots \stackrel{\ell}{\geq} \mathcal{O} \dots \stackrel{L}{\geq} \mathcal{O} \mathbf{y} \right]_i =: \nabla_{ij}^1 \quad (13)$$

$$\begin{cases} \nabla_{ij}^1 &:= \sum_{k_1} \left[\stackrel{1}{\geq}_{ik_1} \mathcal{O}'(y_{k_1}^2) \nabla_{ij}^2 \right] \\ \nabla_{ij}^\ell &:= \sum_{k_\ell} \left[\stackrel{\ell}{\geq}_{k_{\ell-1}k_\ell} \mathcal{O}'(y_{k_\ell}^{\ell+1}) \nabla_{ij}^{\ell+1} \right] \\ \nabla_{ij}^L &:= \stackrel{L}{\geq}_{k_{L-1}j} \mathcal{O}'(y_j) \end{cases} \quad (14)$$

这情况完全类似於经典 Back Prop，以上只是 chain rule 的应用， ∇^ℓ 将每层用 chain rule 分拆开来，所以 ∇ 又叫“local gradient”。上式就是整个网络的反向传递，其中每个 weight 出现 exactly 一次。

但工作还未完，我们要计算 $\frac{\partial J_{ij}}{\partial \geq} = \dot{\nabla}_{ij}^1$ 。（定义 $\geq := \stackrel{\ell}{\geq}_{gh}$ ， $k_0 := i$ ， $k_L := j$ ）

注意： $\mathbf{x} = \mathbf{F}^{-1}(\mathbf{y})$ ，所以 \mathbf{y} 是自变量， \geq 不影响 \mathbf{y} ，所以 $\frac{\partial \mathbf{y}}{\partial \geq} \equiv 0$ 。

上面这句好像有问题：搞不清是 who depends on whom，以下可能错误。

\geq 必会是 $\stackrel{\ell}{\geq}$ 的其中一元，但如果 $\geq \notin \stackrel{\ell}{\geq}$ ，以下的项微分后都会变成 0：

$$\begin{cases} \dot{\nabla}_{ij}^1 = \sum_{k_1} \left[\stackrel{1}{\geq}_{ik_1} \mathcal{O}'(y_{k_1}^2) \dot{\nabla}_{ij}^2 \right] \\ \dot{\nabla}_{ij}^\ell = \sum_{k_\ell} \left[\stackrel{\ell}{\geq}_{k_{\ell-1}k_\ell} \mathcal{O}'(y_{k_\ell}^{\ell+1}) \dot{\nabla}_{ij}^{\ell+1} \right] \\ \dot{\nabla}_{ij}^L = \stackrel{L}{\geq}_{k_{L-1}j} \mathcal{O}'(y_j) \equiv 0 \end{cases} \quad (15)$$

所以实际上只剩下一项：

$$\begin{aligned} \frac{\partial J_{ij}}{\partial \geq} &= \sum_{k_1} \left[\stackrel{1}{\geq}_{k_0k_1} \mathcal{O}'(y_{k_1}^2) \dots \sum_{k_{\ell-2}} \left[\stackrel{\ell-2}{\geq}_{k_{\ell-3}k_{\ell-2}} \mathcal{O}'(y_{k_{\ell-2}}^{\ell-1}) \dots \right. \right. \\ &\quad \left. \left. \begin{cases} \dots \stackrel{\ell-1}{\geq}_{k_{\ell-2},g} \mathcal{O}'(y_g^\ell) \mathcal{O}'(y_h^{\ell+1}) \nabla_{ij}^{\ell+1} \end{cases} \right] \right] \\ &\quad \left. \left[\dots \stackrel{\ell-1}{\geq}_{k_{\ell-2},g} \mathcal{O}'(y_g^\ell) \mathcal{O}'(y_h^{\ell+1}) \right] \right] \quad \text{if } \geq \in \text{last layer} \end{aligned} \quad (16)$$

上式的意思是：每层 layer 重复一块 $[\sum \geq \mathcal{O}']$ ，直到遇到 $\geq = \stackrel{\ell}{\geq}_{gh}$ ，则用结尾形式取代之。

和经典 Back Prop 不同的是，上式只是 $n \times n$ 矩阵中的一个元素，从复杂度而论，每个 weight 的 ∇ 计算，增加了起码 n^2 倍的复杂度（虽然其计算上可以共用一些结果）。记住 $n = \dim$ 状态空间。

可以这样理解：每个 weight 的调教，需要计算这个 weight 对 Jacobian 的影响，而那 Jacobian 是**整个网络**的特性。关键似乎就在於每个 weight 对 Jacobian 的**影响**。

现在回看更高层次的这个式子：

$$\frac{\partial}{\partial w} |J| = \text{tr}(|J| \cdot J^{-1} \cdot \left[\frac{\partial J}{\partial w} \right]) \quad (17)$$

$$= |J| \cdot \text{tr} \left(\left[\frac{\partial y}{\partial x} \right] \cdot \left[\frac{\partial}{\partial w} \frac{\partial x}{\partial y} \right] \right) \quad (18)$$

$$= |J| \cdot \sum_{ij} \left(\frac{\partial y}{\partial x} \right)_{ij} \left(\frac{\partial}{\partial w} \frac{\partial x}{\partial y} \right)_{ij} \quad (19)$$

上式中最重要（最慢）的是那 $(i, j) \in n \times n$ 求和。裡面的第一个因子是 Jacobian J ，第二个因子是我们刚计算了的 $\nabla_w J^{-1}$ 。

Back Prop 的 ∇ 形式上是 $\frac{\partial \text{输出}}{\partial \mathbf{w}}$ ，我们的 ∇ 形式是 $\left[\frac{\partial}{\partial w} \frac{\partial \text{输入}}{\partial \text{输出}} \right]_{n \times n}$ 。

其实我们只需要计算 $\nabla_w |J|$ 的**大约方向**。暂时我在代码中的做法是：忽略式 (19) 中较小的项，那就不需做足 n^2 个乘积。

或者可不可以将 $|J|(\mathbf{w})$ 看成是一个 weight \mathbf{w} 的函数，然后用它的 Taylor series expansion 来近似？

其实更大的问题是：给定一个输入 x ，整个 F 只能输出一个值，因而不能做到 stochastic policy。暂时搁置。

Notes

- Parameters are organized hierarchically (“deep”)
- Jacobian of F at x in terms of W is easy to calculate
-

References